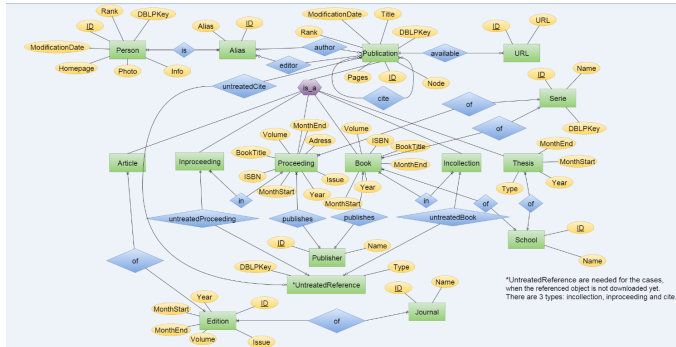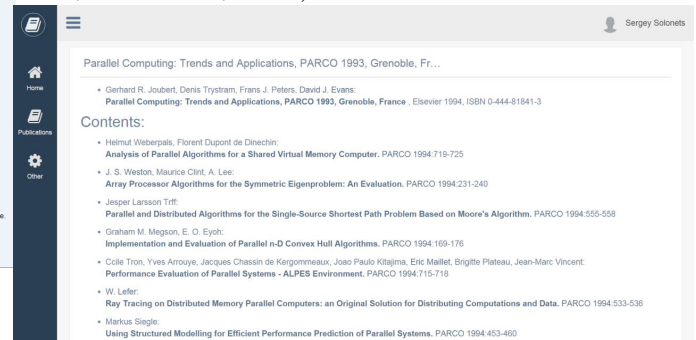# Project report

Sergey Solonets
Innopolis University
420500 Innopolis
Innopolis, Russia
ssolonets@gmail.com

Kozar Anastasia
Innopolis University
420500 Innopolis
Innopolis, Russia
nastya.shev9594@gmail.com

## 1. ER-DIAGRAM



## 2. SITE DESCRIPTION

Structure of our site:

## 2.1 Home

### 2.1.1 Dashboard

Some statistical information.



### 2.1.2 Global search

Allows to search through all types of publications, authors and journals.



## 2.2 Publications

.

Allows to search and view all information about all particular publication types(articles, inproceedings, proceedings, books, incollections, theses)
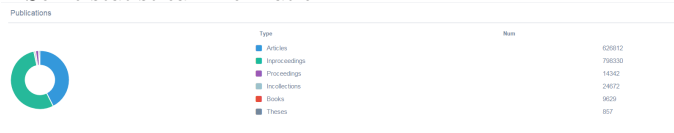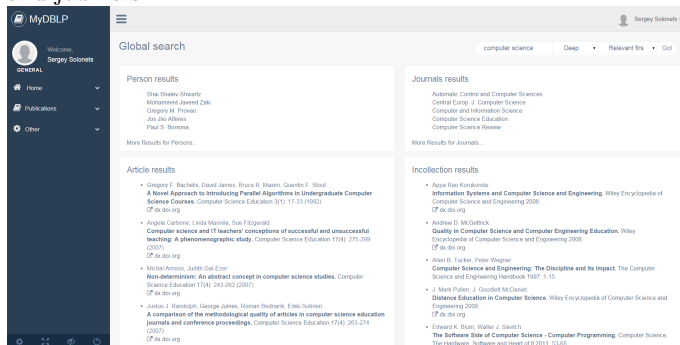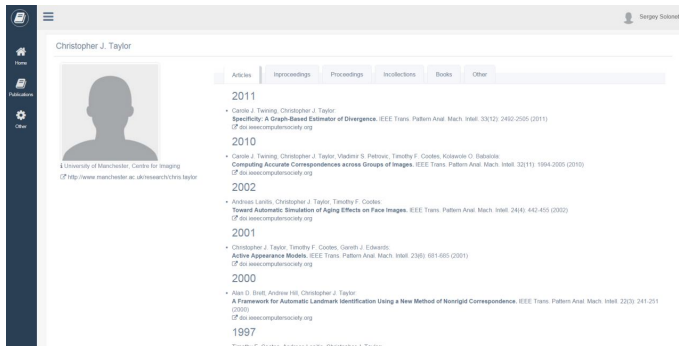


## 2.3 Other

### 2.3.1 Journals

Allows to navigate through journals and editions, which they consist from



### 2.3.2 Persons

Allows to search and view information about persons and their activities

## 3. IMPLEMENTATION DESCRIPTION

All implementation was done on **php** and **MySQL**, using framework **CodeIgniter**. As source was taken **DBLP**. Usage of xml file parsing and DBLP crawling together helped avoid a lot of uncoordinated data in xml file. Operations: Insertion, Update, Deletion implemented in crawler.We had to refuse from their realisation in web interface, because user actions in this case conflict with real-time crawler.
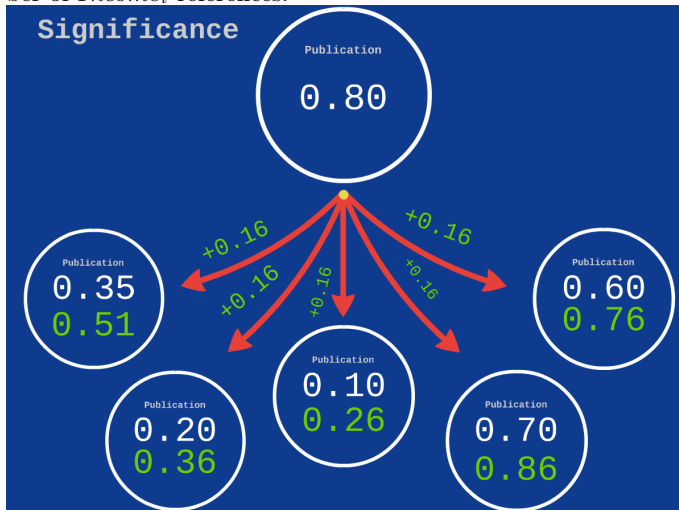
### 3.1 Sorting

Were implemented two sorting methods:

#### 3.1.1 Reasonable first

Sorting publications using rank, which counting as a probability of random surfer to visit this publication by chaotically clicking on references on this site. For this purpose each person and publication getting initial rank and share it proportionally to all correlated entities(person or publication). This is a emulation of Google page rank algorithm. $Rank(Pub) = 0.5 + 0.85 * \sum_{i=1}^{n} \frac{Rank(Income_i)}{C(Income_i)}$, where $Income_i$ - i-th entity referred to this publication, $C(Income_i)$ - number of $Income_i$ references.



#### 3.1.2 Relevant first

Each publication property getting their own weight. Query parse into words. Each searching publication getting weight as sum of weights(multiplied on word length) of attributes, in which these words were found. Also it is used raising co-

efficient in case, when word was found entirely, but not as a word part.



### 3.2 Searching

Were implemented two searching methods Quick and Deep. In both cases query is parsing into words. These word are searched in the specified columns. The difference is only about the columns:

#### 3.2.1 Quick search

The set of columns was selected in such way, that all where clauses would push out all the heavy joins. Speed difference with Deep search in search performance using the same key word can reach 10 times.

#### 3.2.2 Deep search

However Deep search is more convenient, because it searches through all publications information(even those are not displayed in description of this publication). For example, query "Innopolis" would display also publications, which were written in Innopolis.

## 4. DEPLOYMENT

Phase 2 directory contains 2 directories: first one is applications, which is ready to be installed. Just copy it to your root server and edit configuration in **Application/Config**. Second one is routines contains parser, crawler and ranker(rank-crawler) See 2.1.1. MySql structure.sql also there. You can export data from http://vh229.spaceweb.ru/phpMyAdmin using login:solonets password:dblpdblpdblp
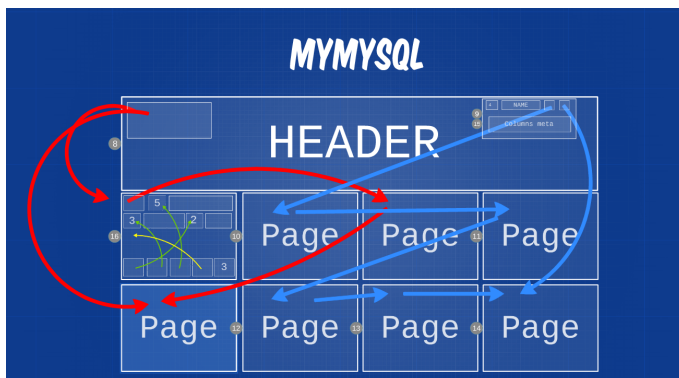
## 5. FILE SYSTEM

We stored all the data in a binary file. The first file page contains meta data, which is consists from:
4bytes(page size) | 4bytes(number of pages) | 4bytes(number of tables) | 4bytes[i-th table name length] | 4bytes[i-th table name] | 4bytes[i-th table start page] | 4bytes[i-th table current page] | 4bytes[i-th table last auto incemented] | 4bytes[i-th table number of attributes] | 4bytes[i-th table i-th attribute name length] | 4bytes[i-th table i-th attribute name] | 4bytes[i-th table i-th attribute type] | byte[i-th table i-th attribute auto incremented] | byte[i-th table i-th attribute indexed]. Let's consider each more closer.
page size - how much bytes was allocated per page

number of pages - how much pages contains this file
number of tables - how much tables contains this file
i-th table name length - name length of the i-th table
i-th table name - name of the i-th table
i-th table start page - from which page starts the records of i-th table
i-th table current page - on which page we will write next time
i-th table last auto incremented - last number used in auto incremented field of i-th table
i-th table number of attributes - number of attributes in i-th table
i-th table i-th attribute name length - i-th table i-th attribute name length
i-th table i-th attribute name - i-th table i-th attribute name
i-th table i-th attribute type - id of an i-th table i-th attribute type(it's int)
i-th table i-th attribute auto incremented - true/false
i-th table i-th attribute indexed - true/false
Each table records stored in their own pages. Page at the beginning contain pointer to the next page. At the end each page contains meta data about the page (pointers to records, last free space, free space quantity).



Indexing implemented by B+ tree. Also we havn't SQL parser, because in this case we would have to implement query optimisation. Using our implementation we can optimise queries by our selves.

```
                SELECT
   * FROM Person JOIN Alias ON Alias.PersonID
         = Person.ID WHERE Person.ID = 20




  db.getTable("Person").selectIndexed(
     "ID",
     new RAInteger(20)
     ).join(
        db.getTable("Alias"),
        new ColumnExpression("ID"),
        "PersonID",
        new EmptyCondition()
        )
     );
```